

Estimation Non-Paramétrique de Densités

Machine Learning

Cédric RICHARD

Université Côte d'Azur

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Objectifs

Nous avons supposé que les fonctions de densité $p(\mathbf{x}|\omega_i)$ sont connues :

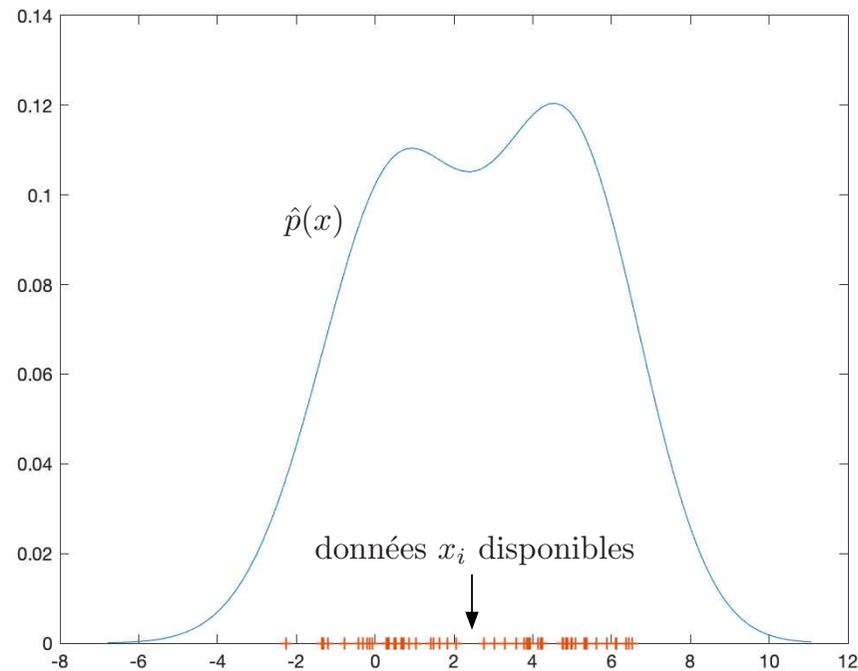
- ▷ soit parfaitement, par leur expression et les paramètres associés
exemple : loi normale $\mathcal{N}(\mu, \sigma^2)$ de moyenne μ et de variance σ^2 connues
- ▷ soit partiellement par leur expression, de paramètres inconnus à estimer
exemple : loi normale $\mathcal{N}(\mu, \sigma^2)$ de moyenne μ et de variance σ^2 à estimer

On ne connaît généralement pas les fonctions de densité de probabilité. Il est alors nécessaire de les inférer à partir des données disponibles

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Exemple (1D)

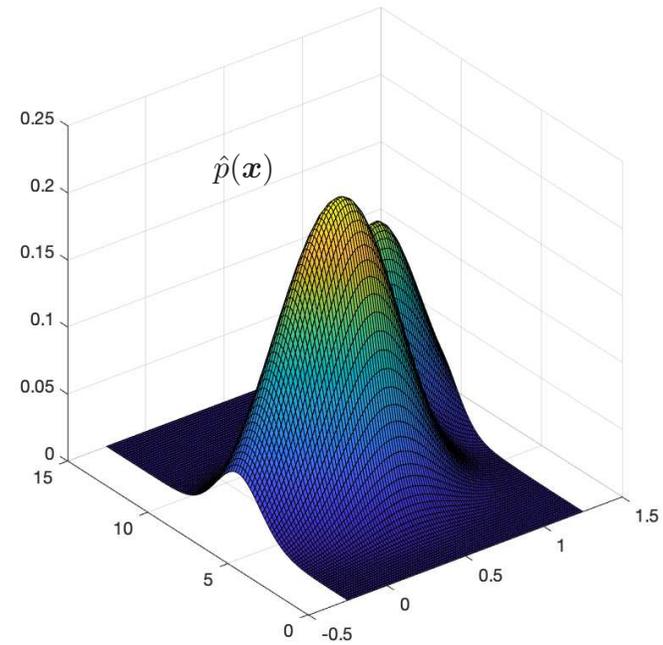
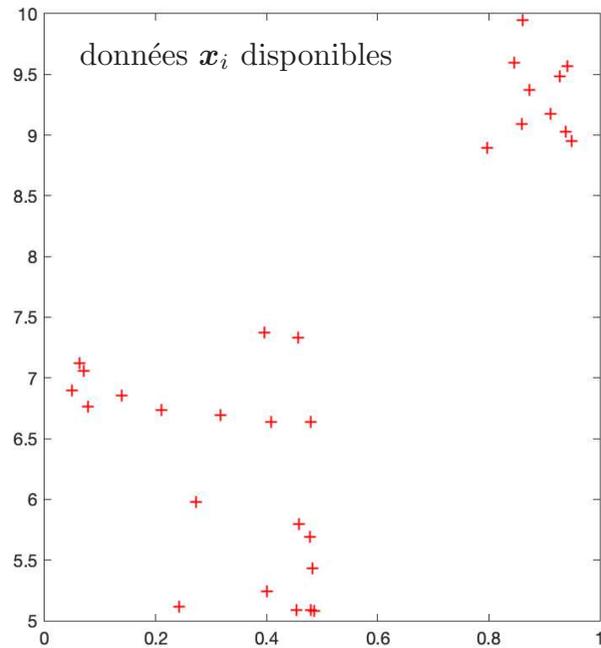
Estimation non-paramétrique de densité attendue :



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Exemple (2D)

Estimation non-paramétrique de densité attendue :



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Histogramme

Une méthode élémentaire est l'histogramme :

- ▷ on subdivise l'espace des observations en boîtes
- ▷ on approxime la densité au centre de chaque boîte par la fraction du nombre de données d'apprentissage en chaque boîte

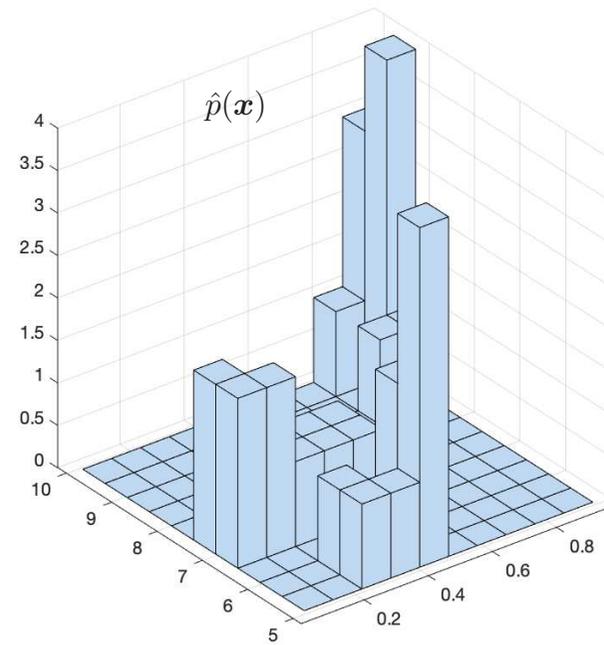
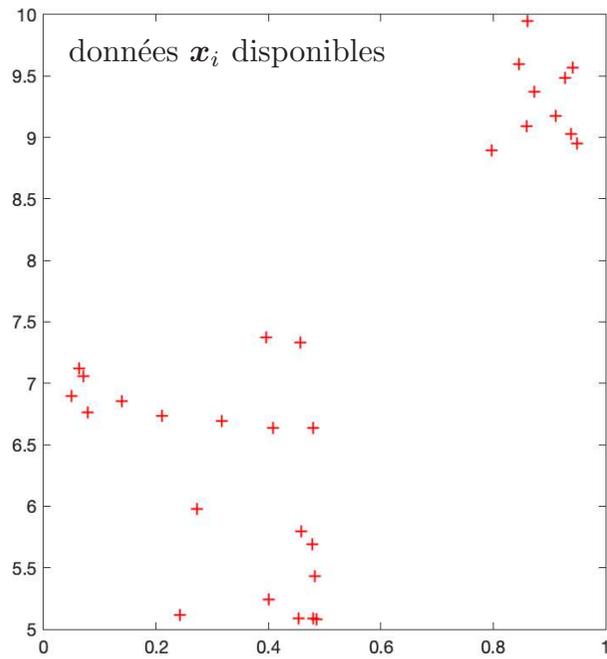
$$\hat{p}(\mathbf{x}) = \frac{1}{N} \frac{\# \text{ de données } \mathbf{x}_k \text{ dans la même boîte que } \mathbf{x}}{\text{volume de la boîte}}$$

Deux paramètres nécessitent d'être fixés :

- ▷ position de la première boîte
- ▷ dimensions des boîtes

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Exemple d'histogramme

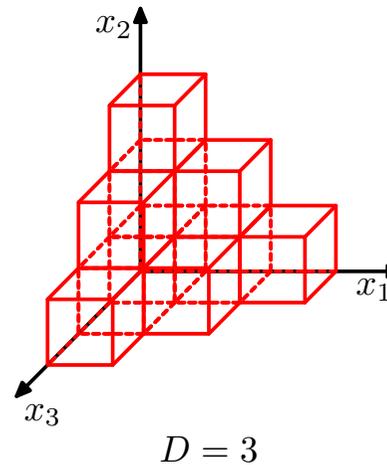
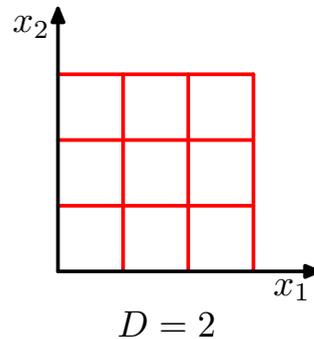
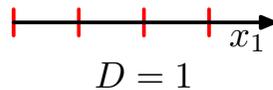


ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Histogramme

Les histogrammes ne sont pas des solutions satisfaisantes, hormis pour une visualisation 1D ou 2D simple :

- ▷ la densité estimée dépend de la position des boîtes et de leur orientation
- ▷ la densité estimée n'est pas lisse
- ▷ la malédiction de la dimensionnalité, due au nombre exponentiel de boîtes nécessaires en fonction de la dimension des données



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Principes généraux

- ▷ La probabilité que \mathbf{x} , distribué selon une loi $p(\mathbf{x})$, appartient à une région \mathcal{R} de l'espace des observations est :

$$P = \int_{\mathbf{x} \in \mathcal{R}} p(\mathbf{x}) d\mathbf{x}$$

- ▷ Soit N observations $\mathbf{x}_1, \dots, \mathbf{x}_N$ distribuées selon $p(\mathbf{x})$. La probabilité que k d'entre elles appartiennent à \mathcal{R} est donnée par la distribution binomiale :

$$P(k) = C_N^k P^k (1 - P)^{N-k}$$

Rappels :

$$E[k] = NP \quad \text{var}[k] = NP(1 - P)$$

- ▷ En conséquence, on a :

$$E\left[\frac{k}{N}\right] = P \quad \text{var}\left[\frac{k}{N}\right] = \frac{P(1 - P)}{N}$$

Le ratio des données dans \mathcal{R} est donc une bonne estimée de P : $P \approx \frac{k}{N}$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Principes généraux

- ▷ Le ratio des données dans \mathcal{R} est une bonne estimée de P : $P \approx \frac{k}{N}$
- ▷ En supposant que \mathcal{R} est suffisamment petite pour que $p(\mathbf{x})$ n'y varie pas de façon significative, on peut écrire

$$P = \int_{\mathbf{x}' \in \mathcal{R}} p(\mathbf{x}') d\mathbf{x}' \approx p(\mathbf{x})V$$

où V est le volume de \mathcal{R}

- ▷ En combinant ces 2 résultats, on aboutit à

$$\boxed{p(\mathbf{x}) \approx \frac{k}{NV}}$$

Cette estimée est d'autant meilleure que N est grand et que V est petit.
En pratique, N est fixé et il faut trouver un compromis entre N et V .

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Principes généraux

Les méthodes d'estimation non-paramétriques de densité reposent sur :

$$p(\mathbf{x}) \approx \frac{k}{NV}$$

Deux stratégies coexistent :

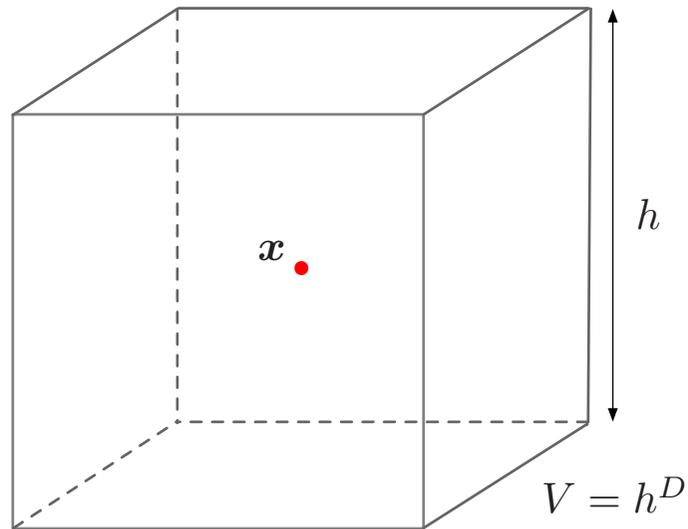
- ▷ on fixe V et on détermine k : méthodes d'estimation à noyau
- ▷ on fixe k et on détermine V : méthode des k -plus-proches-voisins (kPPV)

Ces 2 stratégies font successivement l'objet de ce cours

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Fenêtre de Parzen

On suppose que la région \mathcal{R} enfermant les données est un hypercube dans \mathbb{R}^D , de côté de longueur h , centré en \boldsymbol{x} .



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Fenêtre de Parzen

- ▷ Afin de déterminer le nombre de données contenues dans le cube \mathcal{R} , on introduit la fonction noyau

$$K(\mathbf{u}) = \begin{cases} 1 & |u_j| < 1/2 \quad \forall j = 1, \dots, D \\ 0 & \text{sinon} \end{cases}$$

Il s'agit d'un hypercube de longueur de côté 1 centré sur l'origine

- ▷ On a donc

$$K\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right) = \begin{cases} 1 & \text{si } \mathbf{x}_n \in \mathcal{R} \\ 0 & \text{sinon} \end{cases}$$

- ▷ Le nombre total k de données situées dans \mathcal{R} est

$$k = \sum_{n=1}^N K\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimateur de Parzen

L'estimateur de Parzen au point \boldsymbol{x} est défini par

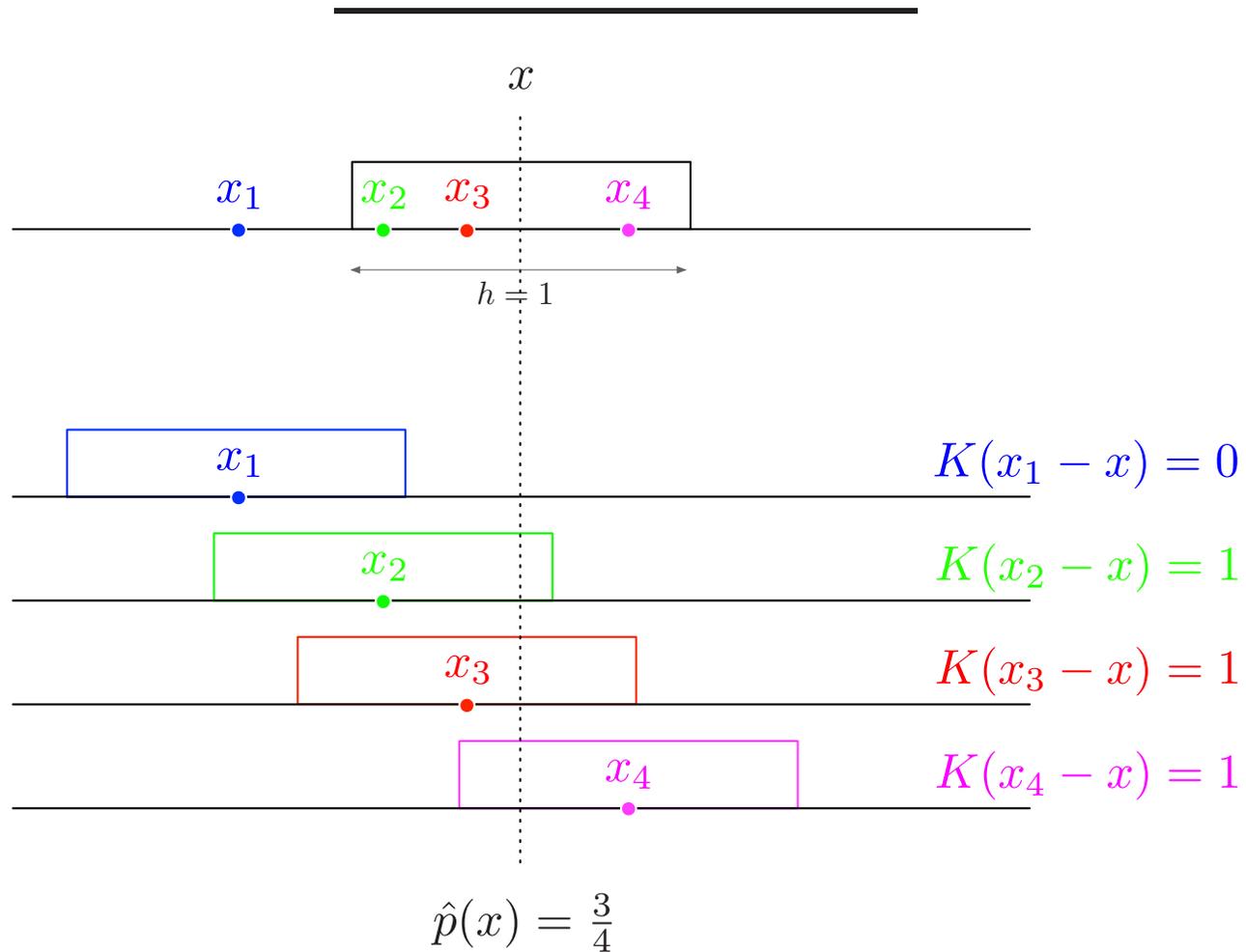
$$\hat{p}(\boldsymbol{x}) = \frac{1}{Nh^D} \sum_{n=1}^N K\left(\frac{\boldsymbol{x} - \boldsymbol{x}_n}{h}\right)$$

Remarque :

Contrairement aux histogrammes, les boîtes de l'estimateur de Parzen sont centrées sur les données

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimateur de Parzen



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimateur de Parzen

On calcule l'espérance mathématique de $\hat{p}(\mathbf{x})$:

$$\begin{aligned} E[\hat{p}(\mathbf{x})] &= \frac{1}{Nh^D} \sum_{n=1}^N E \left[K \left(\frac{\mathbf{x} - \mathbf{x}_n}{h} \right) \right] \\ &= \frac{1}{h^D} E \left[K \left(\frac{\mathbf{x} - \mathbf{x}_n}{h} \right) \right] \\ &= \frac{1}{h^D} \int K \left(\frac{\mathbf{x} - \mathbf{x}'}{h} \right) p(\mathbf{x}') d\mathbf{x}' \end{aligned}$$

où les données \mathbf{x}_n sont supposées indépendantes et distribuées selon $p(\mathbf{x})$.

En conséquence :

$$E[\hat{p}(\mathbf{x})] = p(\mathbf{x}) * K(\mathbf{x}/h)$$

où $*$ est le produit de convolution

Pour $K(\mathbf{x}/h) = \delta(\mathbf{x})$, on a $E[\hat{p}(\mathbf{x})] = p(\mathbf{x})$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimateurs à noyau

La fenêtre rectangulaire a plusieurs inconvénients :

- ▷ la densité estimée a des discontinuités
- ▷ toutes les données \mathbf{x}_i dans le cube ont le même poids, quelle que soit leur distance au centre \mathbf{x}

Afin de pallier ces inconvénients, d'autres fenêtres $K(\mathbf{u})$ existent. Elles sont non-négatives $K(\mathbf{u}) \geq 0$, symétriques $K(\mathbf{u}) = K(-\mathbf{u})$ et vérifient :

$$\int_{\mathbb{R}^D} K(\mathbf{u}) d\mathbf{u} = 1$$

Il en résulte l'estimateur :

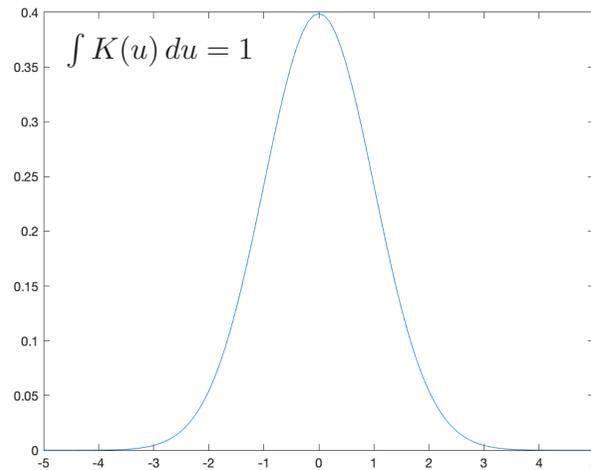
$$\hat{p}(\mathbf{x}) = \frac{1}{Nh^D} \sum_{n=1}^N K\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Exemple d'estimateur à noyau

La fenêtre gaussienne :

$$K(\mathbf{u}) = \frac{1}{(2\pi)^{\frac{D}{2}}} \exp\left(-\frac{\|\mathbf{u}\|^2}{2}\right)$$

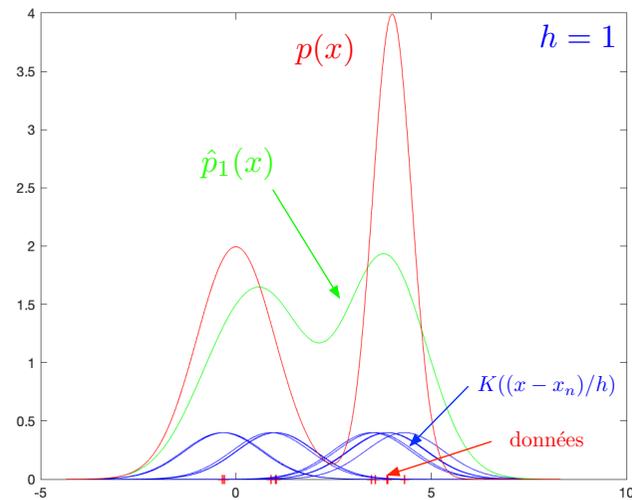


ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimateurs à noyau

$$\hat{p}(x) = \frac{1}{Nh^D} \sum_{n=1}^N K\left(\frac{x - x_n}{h}\right)$$

- ▷ la densité estimée est une somme de motifs centrés sur les données
- ▷ la fonction noyau $K(\mathbf{u})$ détermine la forme des motifs
- ▷ la largeur de bande h , ou paramètre de lissage, définit la largeur des motifs

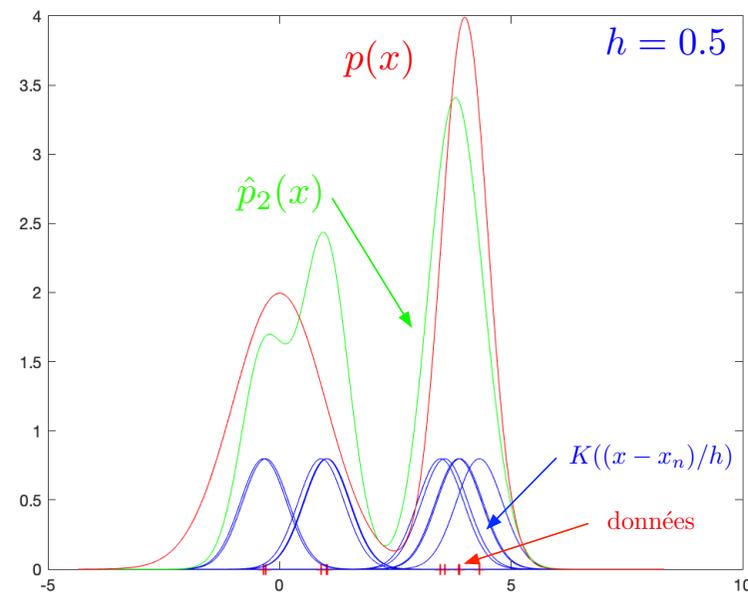
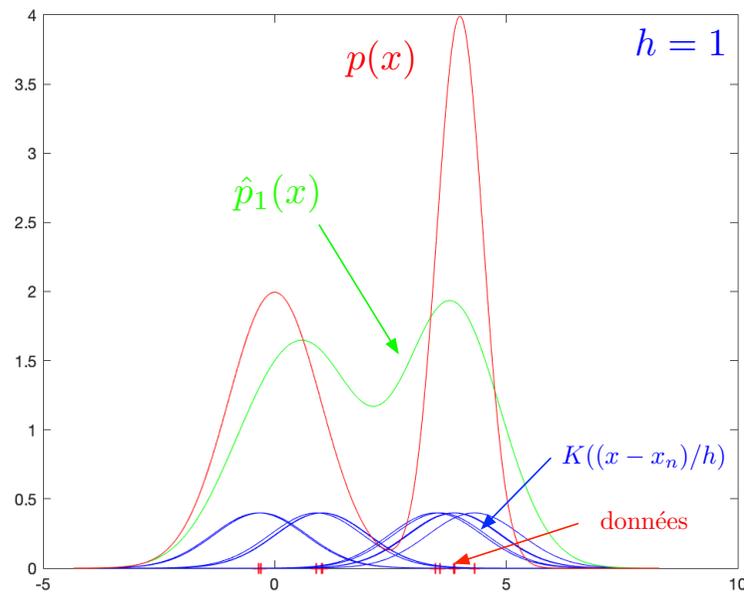


ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Largeur de bande h

Le choix de la largeur de bande h est crucial :

- ▷ $h \nearrow$: densité estimée $\hat{p}(x)$ lisse, masquant éventuellement les modes de $p(x)$
- ▷ $h \searrow$: densité estimée $\hat{p}(x)$ piquée, potentiellement difficile à interpréter



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimation de la largeur de bande h

Le choix subjectif de h n'est pas forcément possible, lorsque la dimension D des données est grande.

Il existe des méthodes automatiques pour le faire, reposant elles-mêmes sur des hypothèses pas nécessairement vérifiées

Exemple :

- ▷ On fait l'hypothèse d'une distribution standard $p(\mathbf{x})$, et on recherche h optimum au sens de l'erreur quadratique moyenne :

$$\hat{h} = \arg \min_h E \left[\int [p(\mathbf{x}) - \hat{p}(\mathbf{x})]^2 d\mathbf{x} \right]$$

- ▷ En choisissant une loi normale pour $p(\mathbf{x})$, on trouve

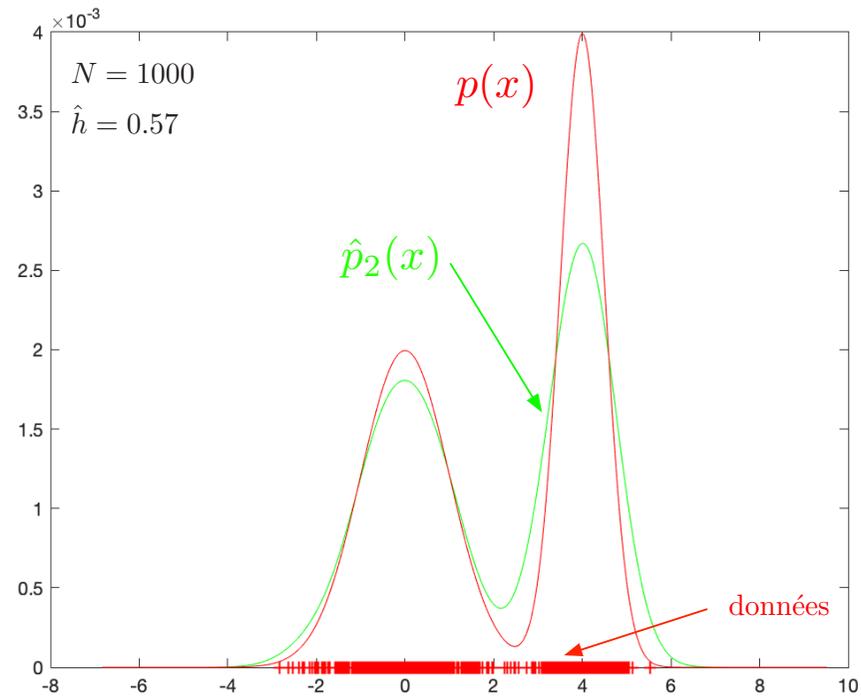
$$\hat{h} = 1.06 \hat{\sigma} N^{-1/5}$$

où $\hat{\sigma}$ est l'écart-type estimé à partir des données disponibles

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimation de la largeur de bande h

Par la méthode précédente ($h^* = 1.06 \hat{\sigma} N^{-1/5}$) :



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimation de densités multi-variées

Les densités multi-variées, i.e., $p(\mathbf{x})$ avec $\mathbf{x} \in \mathbb{R}^D$, peuvent être estimées à l'aide de noyaux scalaires, i.e., $K(u)$ avec $u \in \mathbb{R}$

Méthode : produit de noyaux = noyau produit

$$\hat{p}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N K(\mathbf{x}; \mathbf{x}_n, h_1, \dots, h_D)$$
$$\text{avec } K(\mathbf{x}; \mathbf{x}_n, h_1, \dots, h_D) = \frac{1}{h_1 \dots h_D} \prod_{d=1}^D K_d \left(\frac{x_d - x_{n,d}}{h_d} \right)$$

où

$K_d(u)$ est le noyau scalaire associé à la dimension d

h_d est la largeur de bande associée au noyau $K_d(u)$

$x_{n,d}$ est la d^{e} composante de \mathbf{x}_n

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Estimation de densités multi-variées

L'usage d'un noyau produit n'induit pas que les variables de \mathbf{x} sont indépendantes

Si les variables étaient indépendantes, on pourrait écrire :

$$p(\mathbf{x}) = \prod_{d=1}^D p_d(x_d)$$

que l'on estimerait par :

$$\begin{aligned} \hat{p}(\mathbf{x}) &= \prod_{d=1}^D \hat{p}_d(x_d) \\ &= \prod_{d=1}^D \left(\frac{1}{Nh_d} \sum_{n=1}^N K\left(\frac{x_d - x_{n,d}}{h_d}\right) \right) \end{aligned}$$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

Conclusions sur les estimateurs à noyau

Avantages :

- ▷ applicable à des données issues de n'importe quelle distribution $p(\mathbf{x})$
- ▷ convergence asymptotique de $\hat{p}(\mathbf{x})$ vers $p(\mathbf{x})$ pour $N \rightarrow +\infty$

Inconvénients :

- ▷ nécessite de disposer d'un grand nombre de données
- ▷ choix de la largeur de bande h de la fenêtre non trivial
- ▷ nécessite de disposer de moyens de calcul performants

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k-plus-proches-voisins

Les méthodes d'estimation non-paramétriques de densité reposent sur :

$$p(\mathbf{x}) \approx \frac{k}{NV}$$

où

$$\left\{ \begin{array}{l} k \quad \text{nombre de données voisines de } \mathbf{x} \\ V \quad \text{volume du voisinage de } \mathbf{x} \\ N \quad \text{nombre total de données} \end{array} \right.$$

Deux stratégies coexistent :

- ▷ on fixe V et on détermine k : méthodes d'estimation à noyau
- ▷ on fixe k et on détermine V : méthode des k -plus-proches-voisins (kPPV)

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k-plus-proches-voisins

Principe :

- ▷ On fait croître le volume du voisinage de \boldsymbol{x} jusqu'à ce qu'il inclut k voisins
- ▷ On évalue le volume V du voisinage de \boldsymbol{x}

$$V = c_D r_k^D(\boldsymbol{x}) \text{ avec } c_D = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)}$$

où

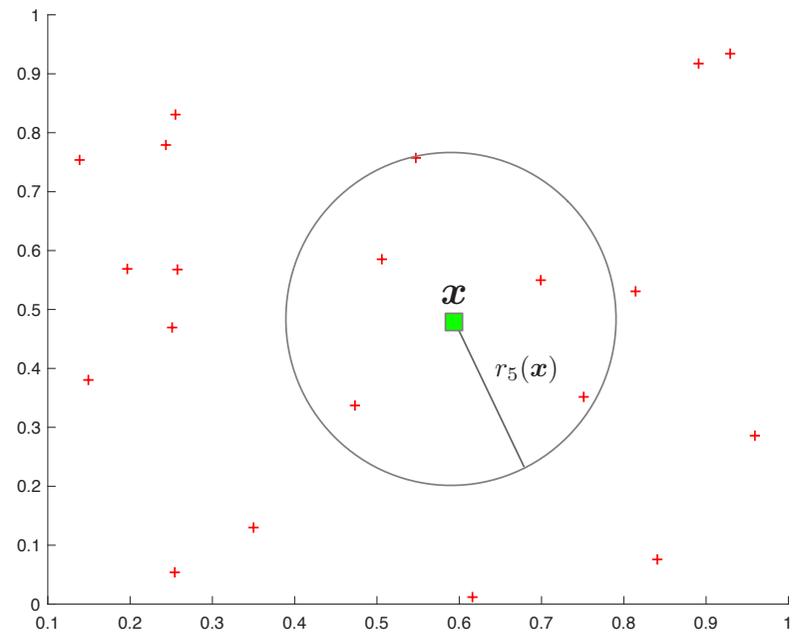
$r_k(\boldsymbol{x})$ est la distance euclidienne entre \boldsymbol{x} et son k^e plus proche voisin

$\Gamma(n + 1) = n\Gamma(n)$ avec $\Gamma(1/2) = \sqrt{\pi}$ et $\Gamma(1) = 1$

D	$\Gamma(\frac{D}{2} + 1)$	c_D	V
1	$\frac{1}{2}\sqrt{\pi}$	2	$2r_k$
2	1	π	πr_k^2
3	$\frac{3}{4}\sqrt{\pi}$	$\frac{4}{3}\pi$	$\frac{4}{3}\pi r_k^3$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k-plus-proches-voisins

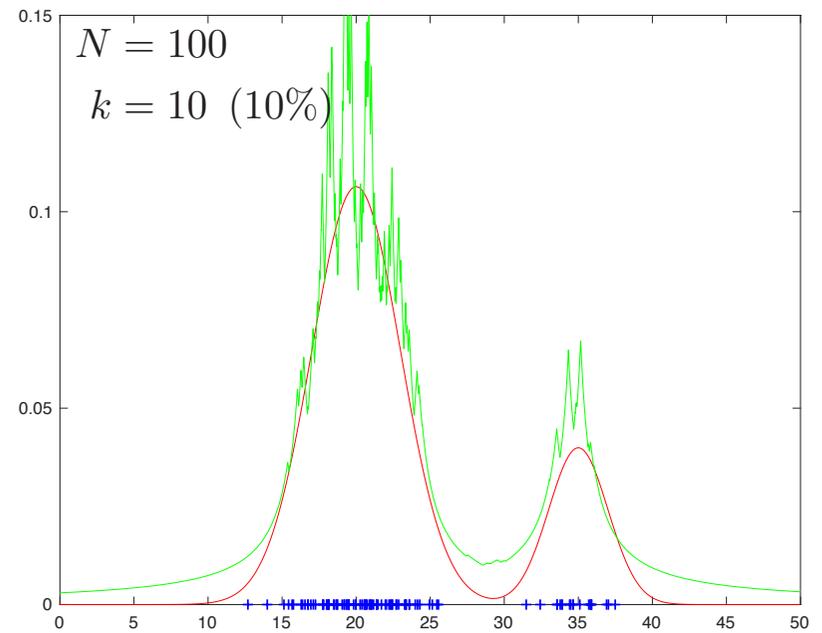
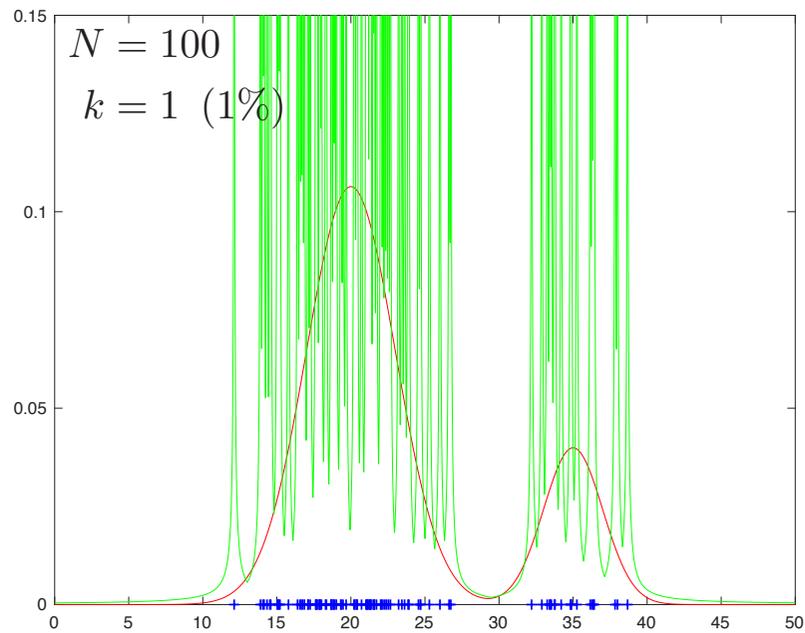
L'estimation de densité par *k*-plus-proches-voisins n'est pas satisfaisante car

- ▷ les estimées sont sujettes au bruit
- ▷ les estimées présentes des queues lourdes
- ▷ les estimées présentes des discontinuités
- ▷ les estimées ne sont pas des densités de probabilité car leur intégrale sur l'ensemble des observations diverge

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : exemple 1D

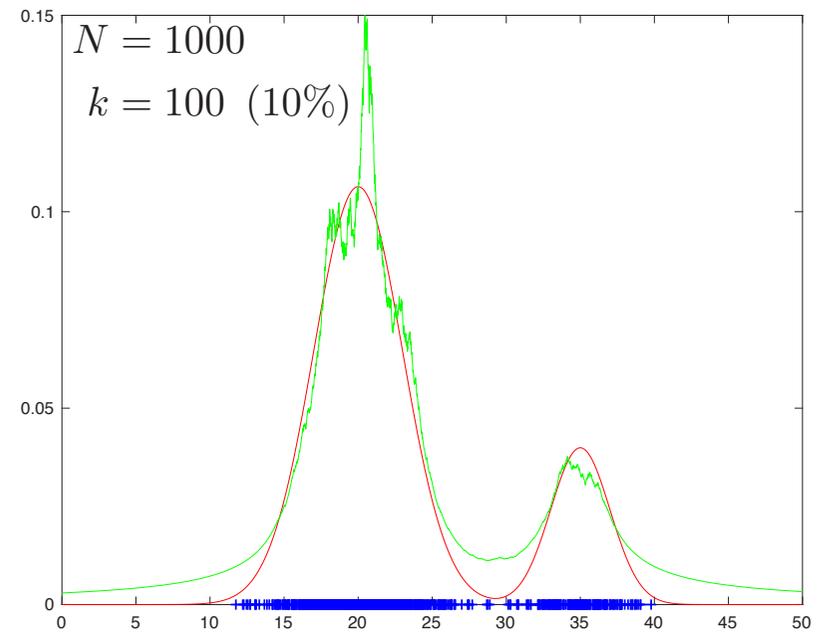
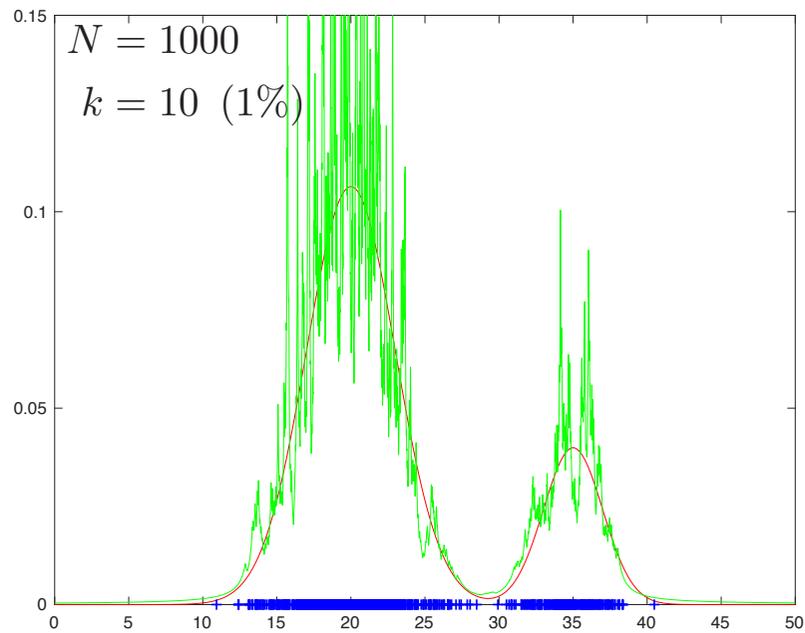
Illustration avec $p(x) = 0.8 \mathcal{N}(20, 9) + 0.2 \mathcal{N}(35, 4)$



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : exemple 1D

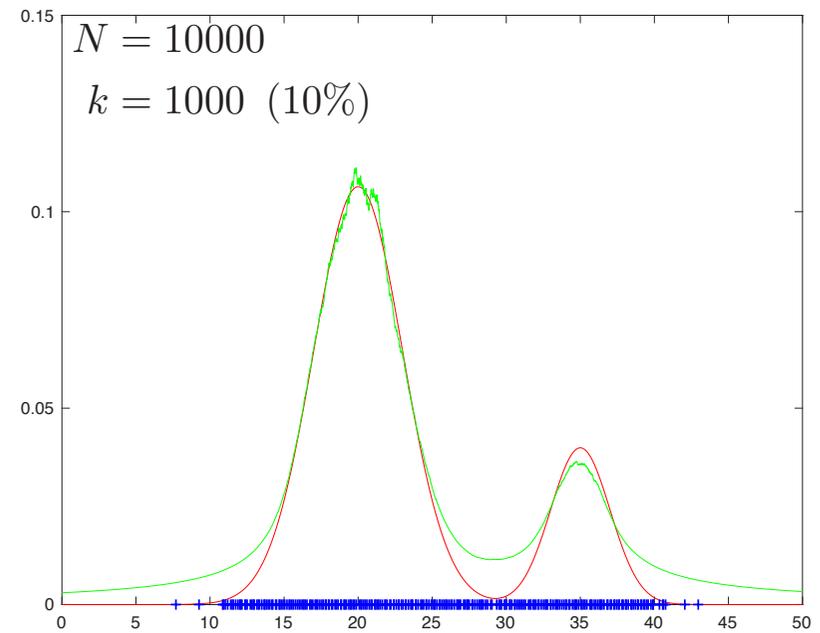
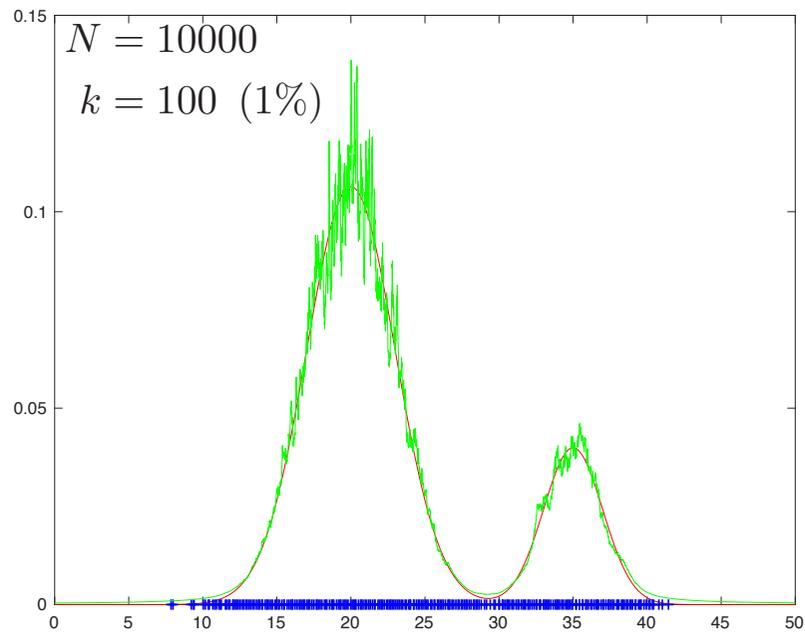
Illustration avec $p(x) = 0.8 \mathcal{N}(20, 9) + 0.2 \mathcal{N}(35, 4)$



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : exemple 1D

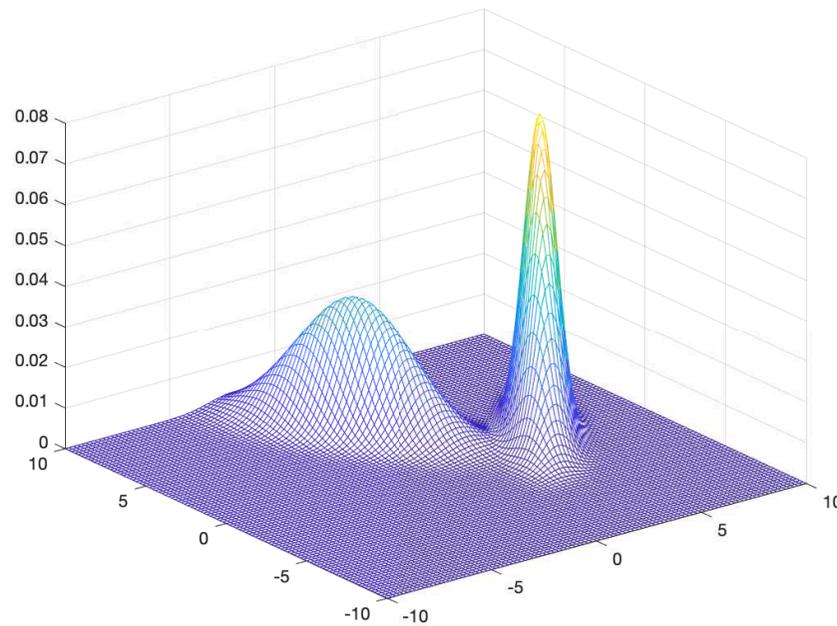
Illustration avec $p(x) = 0.8 \mathcal{N}(20, 9) + 0.2 \mathcal{N}(35, 4)$



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : exemple 2D

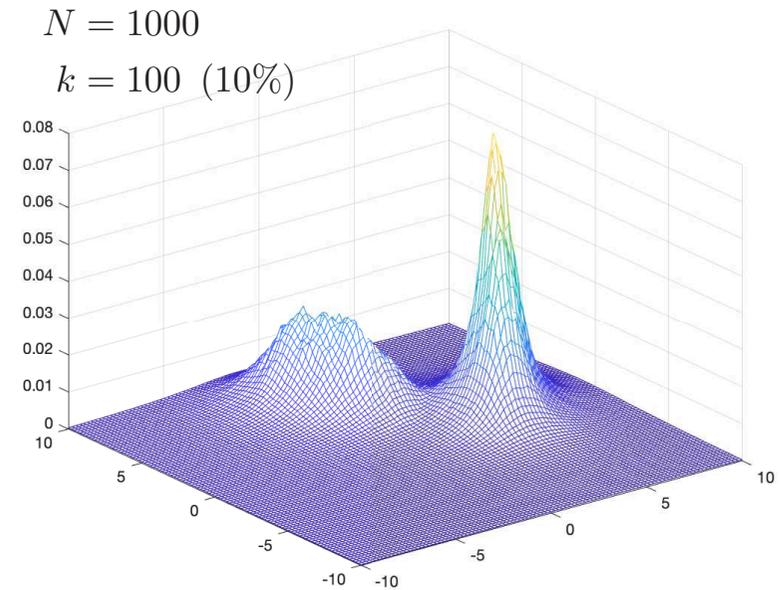
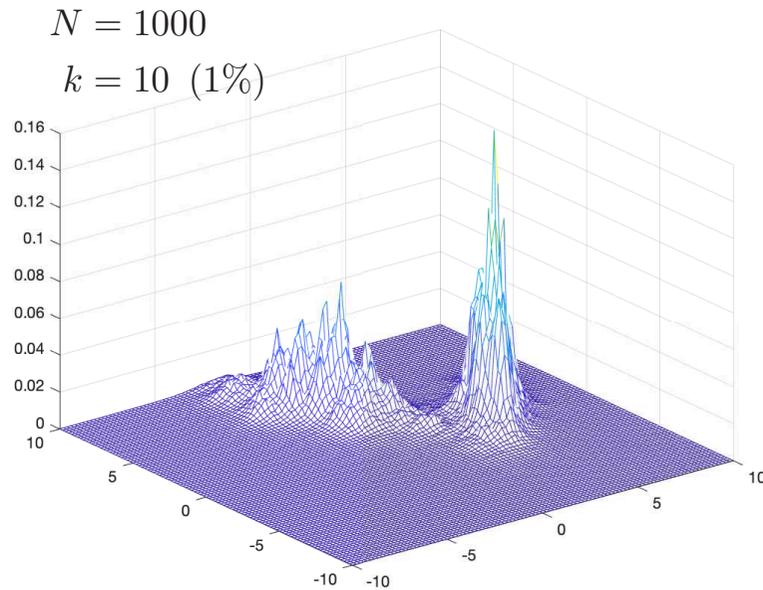
Illustration avec $p(x) = 0.5 \mathcal{N} \left(\begin{bmatrix} 0 \\ 5 \end{bmatrix}, \begin{bmatrix} 2 & -2 \\ -2 & 5 \end{bmatrix} \right) + 0.5 \mathcal{N} \left(\begin{bmatrix} 5 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix} \right)$



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : exemple 2D

Illustration avec $p(x) = 0.5 \mathcal{N} \left(\begin{bmatrix} 0 \\ 5 \end{bmatrix}, \begin{bmatrix} 2 & -2 \\ -2 & 5 \end{bmatrix} \right) + 0.5 \mathcal{N} \left(\begin{bmatrix} 5 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix} \right)$



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k-plus-proches-voisins : application à la classification

L'estimation de densité par *k*-plus-proches-voisins n'est pas satisfaisante. Toutefois, cette méthode offre une approximation simple du classifieur optimum de Bayes.

Problème : N données d'apprentissage, divisées en N_i données par classe ω_i

On estime

- ▷ les fonctions de vraisemblance par : $p(\mathbf{x}|\omega_i) = \frac{k_i}{N_i V}$ avec V le voisinage de \mathbf{x}
- ▷ la fonction de densité de probabilité : $p(\mathbf{x}) = \frac{k}{NV}$
- ▷ les probabilités *a priori* : $P(\omega_i) = \frac{N_i}{N}$

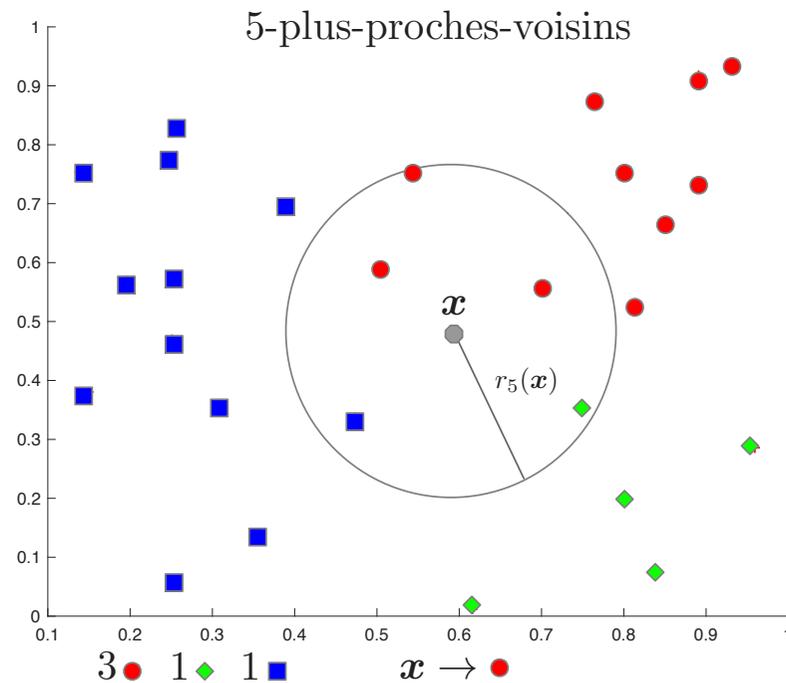
Les fonctions de densité de probabilité *a posteriori* sont données par :

$$p(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})} = \frac{\frac{k_i}{N_i V} \frac{N_i}{N}}{\frac{k}{NV}} = \frac{k_i}{k}$$

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : application à la classification

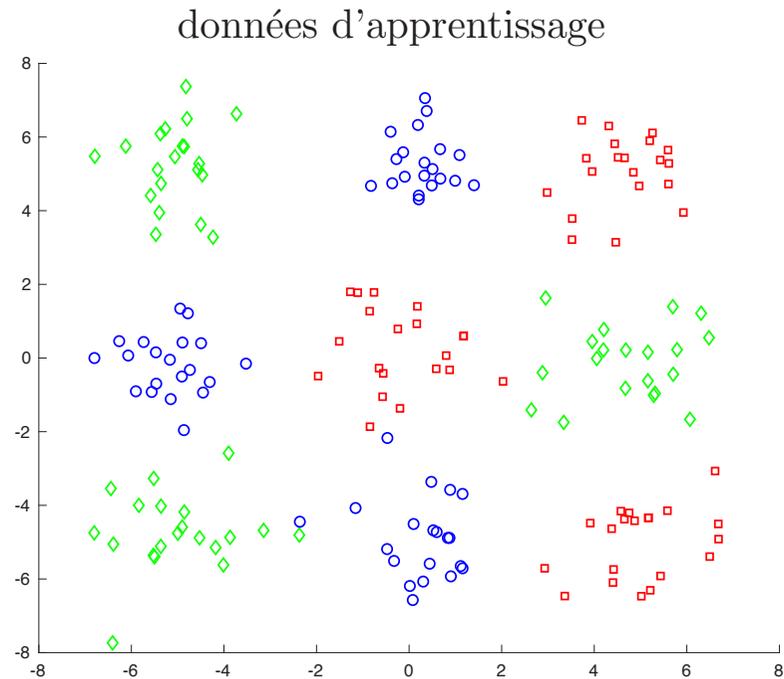
La méthode des k -plus-proches-voisins associe la donnée \boldsymbol{x} à la classe ω_i majoritairement représentée parmi ses k voisins



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : application à la classification

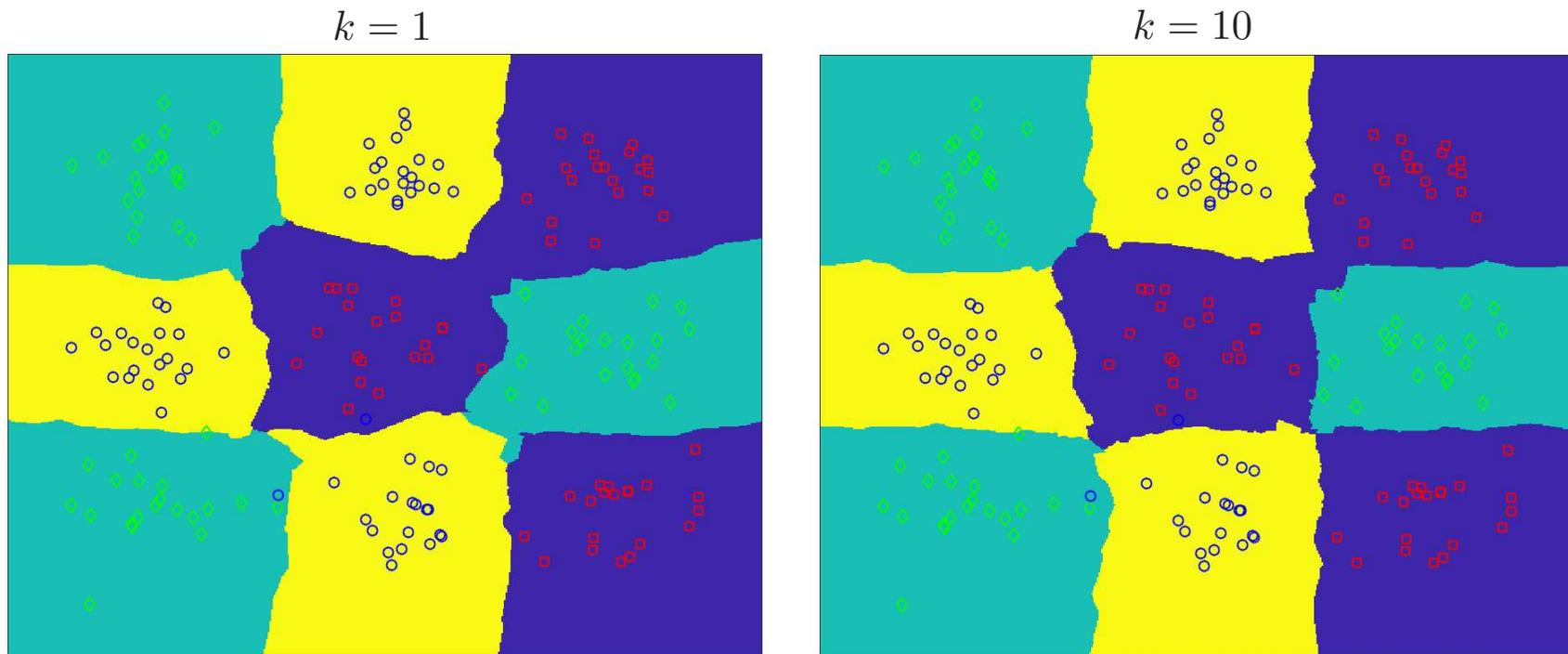
Exemple de classification à 3 classes (rouge, vert, bleu), où chaque cluster de 20 points suit une loi normale de covariance identité.



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : application à la classification

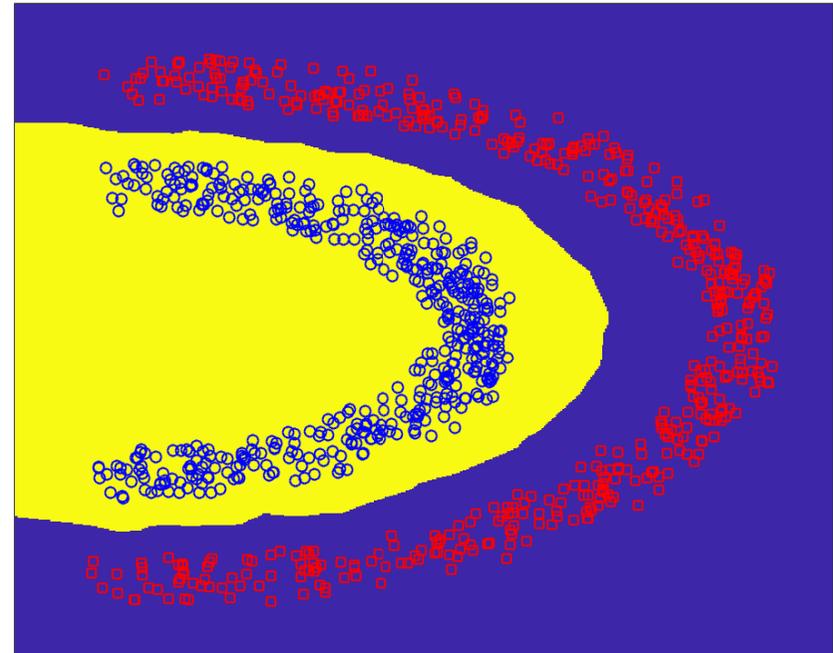
Exemple de classification à 3 classes (rouge, vert, bleu), où chaque cluster de 20 points suit une loi normale de covariance identité.



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : application à la classification

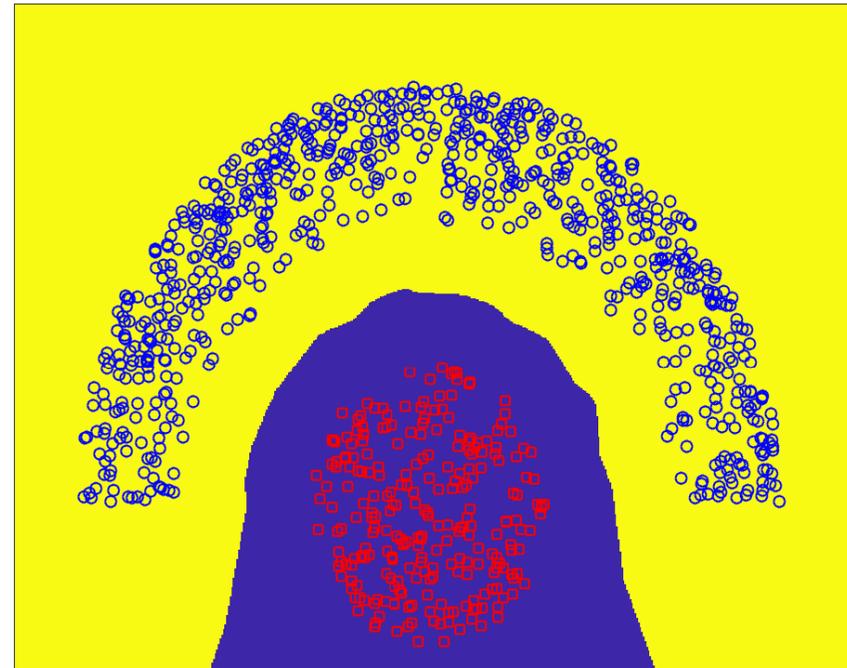
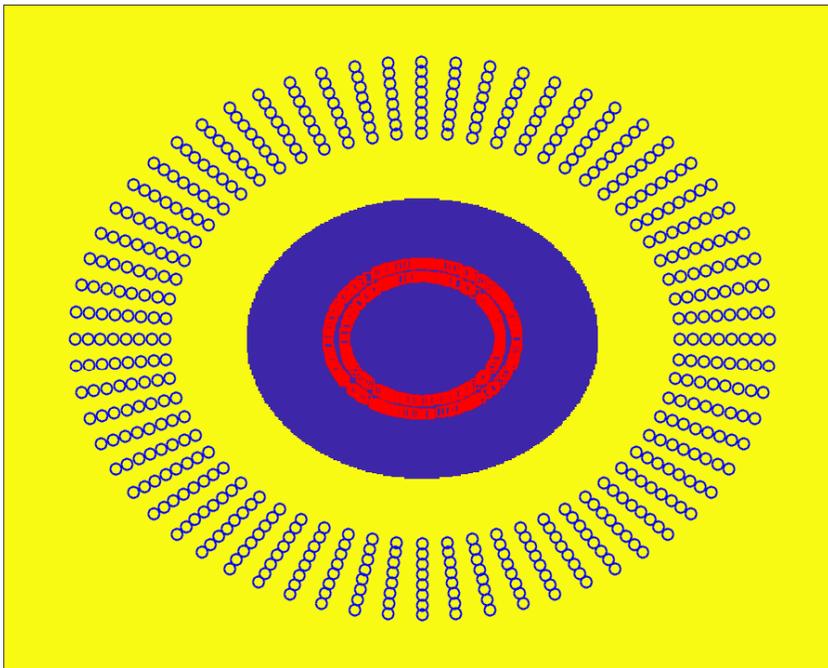
Résultats de la classification ($k = 1$)



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : application à la classification

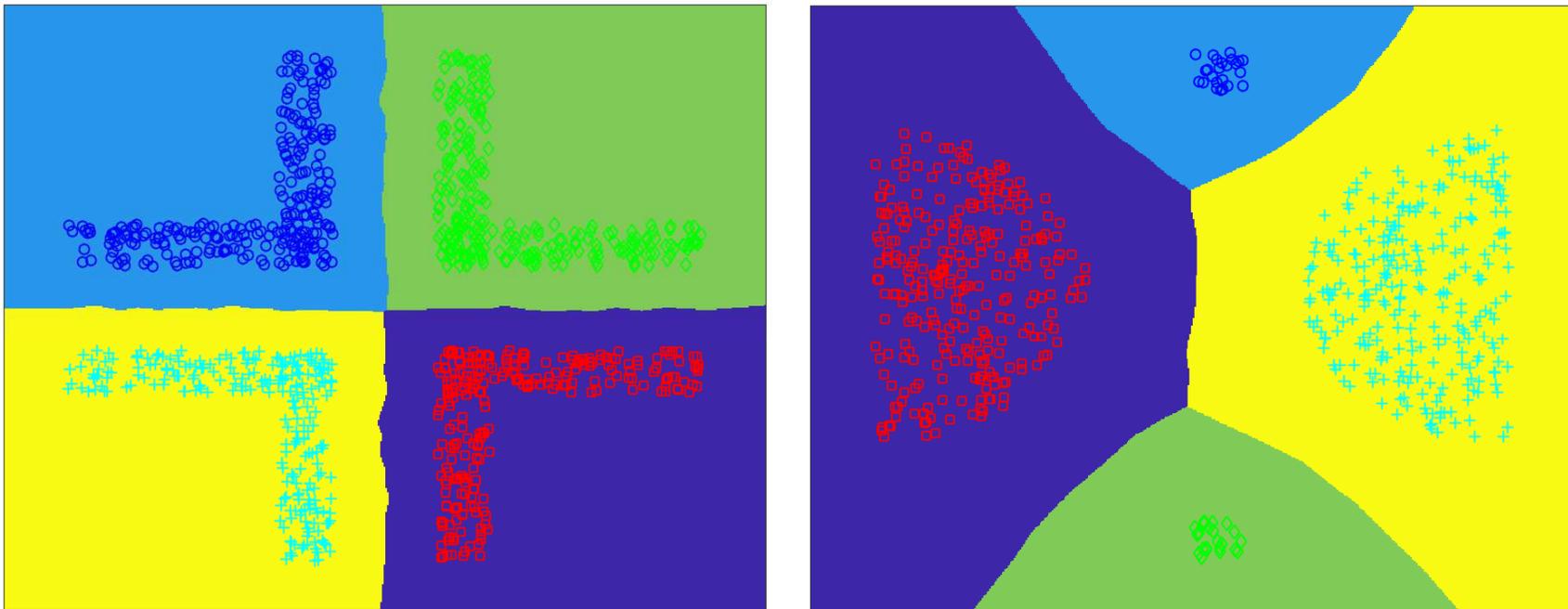
Résultats de la classification ($k = 1$)



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k -plus-proches-voisins : application à la classification

Résultats de la classification ($k = 1$)



ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k-plus-proches-voisins : application à la classification

L'algorithme des *k*-plus-proches-voisins n'extrait pas d'information des données afin d'estimer les paramètres d'un modèle.

En conséquence, la phase d'apprentissage demande pas/peu de ressources contrairement à la mise en œuvre.

Avantages :

- ▷ mise en œuvre simple, et implantation parallèle possible
- ▷ performant puisque, pour $N \rightarrow +\infty$, on a :

$$P_{Bayes}(\text{erreur}) < P_{1PPV}(\text{erreur}) < 2P_{Bayes}(\text{erreur})$$

- ▷ capacité d'adaptation car utilise l'information locale

Inconvénients :

- ▷ nécessite une grande capacité de stockage
- ▷ classification lourde d'un point de vue calculatoire
- ▷ sensible à la malédiction de la dimensionnalité

ESTIMATION NON-PARAMÉTRIQUE DE DENSITÉ

k-plus-proches-voisins : application à la classification

De nombreuses stratégies d'améliorations existent pour

- ▷ réduire la taille de la base d'apprentissage en exploitant sa redondance
 - en supprimant les données mal-classées, donc ambiguës
 - en supprimant les données bien-classées, donc loin de la frontière
- ▷ réduire le temps de recherche des voisins
 - en partitionnant l'espace en cellules (diagramme de Voronoï)
 - en hiérarchisant les données sous forme d'arbre
- ▷ choisir la métrique utilisée pour le calcul de distance